# THEORY HAI ZAROORI NOTES

# SESSION 1

# CA. PRANAV POPAT

## SESSION LINK:
https://www.youtube.com/live/58T5_vLBbI8?si=V-tknpZDGvRO-GX1

## JOIN TELEGRAM CHANNEL FOR ALL UPDATES AND NOTES:
https://telegram.me/learnwithpranav

## THEORY WEIGHTAGE

| Chapter | 13. Statistical Description of Data | 14. Central Tendency & Dispersion | 15. Probability | 16. Theoretical Distribution | 17. Correlation & Regression | 18. Index Numbers | Total |
|---|---|---|---|---|---|---|---|
| May 18 | 2 | 4 | 2 | 3 | 6 | 8 | 25 |
| Nov 18 | 6 | 1 | 0 | 0 | 2 | 3 | 12 |
| Jun 19 | 5 | 3 | 1 | 0 | 1 | 5 | 15 |
| Nov 19 | 1 | 7 | 0 | 2 | 2 | 5 | 17 |
| Nov 20 | 8 | 5 | 0 | 4 | 3 | 6 | 26 |
| Jan 21 | 10 | 5 | 1 | 2 | 2 | 4 | 24 |
| Jul 21 | 6 | 1 | 0 | 0 | 1 | 0 | 8 |
| Dec 21 | 3 | 5 | 0 | 0 | 2 | 4 | 14 |
| Jun 22 | 9 | 3 | 0 | 1 | 4 | 6 | 23 |
| Dec 22 | 4 | 3 | 1 | 2 | 1 | 3 | 14 |
| Jun 23 | 2 | 0 | 0 | 0 | 0 | 2 | 4 |

**THEORY CONCEPTS**

**Statistical Description of Data – Basics of Statistics**

| | |
|---|---|
| **Definition of Statistics** | • Plural Sense: Any data – quantitative or qualitative used for statistical analysis.<br>• Singular Sense: Scientific method of collecting, analyzing, and presenting data to draw statistical inferences. It is also called as Science of Averages or Science of Counting *(Conclusion)* |

| **Origin of Word** | Language | Actual Word | Memorize by |
|---|---|---|---|
| | Latin | Status | Latus |
| | Italian | Statista | Pasta |
| | German | Statistic | Breadstick |
| | French | Statistique | Barbeque |

| | | |
|---|---|---|
| **Publication** | **Koutilya's Arthashastra** | • Record of Birth and Deaths<br>• Chandragupta's reign<br>• 4th Century B.C |
| | **Abu Fezal's Ain-i-Akbari** | • Record on Agriculture<br>• Akbar Reign<br>• 16th Century A.D. |
| | **First Census** | • Egypt 300 BC to 2000 BC<br>• By Pharaoh |

| | |
|---|---|
| **Application of Statistics** | • Economics: Demand Analysis, Future Projection etc.<br>• Business Management: Decision making using quantitative techniques not intuition<br>• Industry and Commerce: Profit maximization using business data – sales, purchase, market etc. by consulting experts |
| **Limitation of Statistics** | • It deals with aggregate data and not individual data<br>• Quantitative data can only be used, however for qualitative – it needs to be converted into quantitative<br>• Projections are based on conditions/ assumptions and any change in that will change the projection. Example: Future projections of sales<br>• Sampling based conclusions are used, improper sampling leads to improper results. Random Sampling is must. |
| **Data** | • Quantitative Information shown as number<br>• Primary: first time collected by agency/ investigator<br>• Secondary: collected data used by different person/ agency |

| | | |
|---|---|---|
| **Variable** | • Measurable Data – Value can vary | |
| | **Discrete Variable** | • When a variable assumes a finite or count ably infinite isolated values.<br>• Example: no. of petals in a flower, no. of road accident in locality |
| | **Continuous Variable** | • When a variable assumes any value from the given interval (can also be in decimals, fractions).<br>• Example: height, weight, sale, money |
| **Attribute** | • Qualitative Characteristics. Example: gender of a baby, the nationality of a person, the colour of a flower etc. | |

| **Collection of Primary Data – Interview Method** | **Method** | **Details** |
|---|---|---|
| | **Personal Interview** | • Where data is collected directly from respondents.<br>• Highly Accurate – Low Coverage<br>• Example: Natural Calamity, Door to Door Survey |
| | **Indirect Interview** | • When reaching respondent is difficult, data is collected by contacting associated persons.<br>• Highly Accurate – Low Coverage<br>• Example: Rail accident |
| | **Telephone Interview** | • Data is collected over phone<br>• Quick and non-expensive method<br>• Low Accuracy – High Coverage |

| | |
|---|---|
| **Collection of Primary Data – Mailed Questionnaire Method** | • In this method well drafted and soundly sequenced questionnaire,<br>• covering all the important aspects of the data requirement is sent to respondent for filling.<br>• Here coverage is wide but amount of non-responses will be maximum |
| **Collection of Primary Data – Observation Method** | • In this method data is collected by direct observation or using instrument.<br>• For example: data on height and weight for a group of students.<br>• Although more accurate but it is time consuming, low coverage and laborious method. |
| **Collection of Primary Data – Questionnaire Filled and sent by Enumerators** | • Mix of Interview and Mailed Questionnaire<br>• Enumerator means a Person who directly interacts with respondent and fills the questionnaire.<br>• It is generally used in case of Surveys and Census. |

| | | |
|---|---|---|
| **Sources of Secondary Data** | **International Sources** | World Health Organization (WHO), International Monetary Fund (IMF), International Labor Organization (ILO), World Bank |
| | **Government Sources** | In India – Central Statistics Office (CSO), Indian Agricultural Statistics by the Ministry of Food and Agri, National Sample Survey Office- NSSO, Regulators – RBI, SEBI, RERA, IRDA |
| | **Private or Quasi-govt. sources** | Indian Statistical Institute (ISI), Indian Council of Agriculture, NCERT |

| | |
|---|---|
| **Scrutiny of Data** | • checking accuracy and consistency of data<br>• There is no rule for it, one must apply his intelligence, patience and experience while scrutinizing the given information.<br>• Internal Consistency: When two or more series of related data are given, we should check consistency among them. |

| | | |
|---|---|---|
| **Presentation of Data – Classification / Organization of Data** | **Classification or Organisation**: putting data in a neat, precise, and condensed form, making it comparable, suitable for analysis, more understandable. | |
| | **Chronological/ Temporal/ Time Series Data** | • Data arranged based on Time<br>• Example: Revenues YoY i.e year on year |
| | **Geographical or Spatial Series Data** | • Arrangement based on regions<br>• Example: Country wise Revenue of a global company |
| | **Qualitative or Ordinal Data** | • Based on some attribute<br>• Nationality Wise Medal Winners in Olympics |
| | **Quantitative or Cardinal Data** | • Based on some variable<br>• Example: Frequency Distribution of a Data |

| | |
|---|---|
| **Mode of Presentation of Data – Textual** | • This method comprises presenting data with the help of a paragraph or several paragraphs.<br>• This is not a suitable mode of presentation as it is dull, monotonous and non-comparable. |

| | |
|---|---|
| **Mode of Presentation of Data – Tabular Form** | • When data is shown in the form of **Table**.<br>• Useful in easy comparison<br>• Complicated data can be presented<br>• Table is must to create a diagram<br>• No analysis possible without table<br>• Components of Table |

| | | |
|---|---|---|
| **Components of Table** | **Description** | **Name of Component of Table** |
| | **Entire Upper Part** | Box Head |
| | **Upper Part describing columns and sub-columns** | Caption |
| | **Left part of the table describing rows** | Stub |
| | **Main Data of Table** | Body |
| | **Source of Data at the bottom of Table** | Footnote |

| | |
|---|---|
| **Mode of Presentation of Data – Diagrams** | • Can be used by educated and uneducated section of society<br>• Hidden trend can be traced<br>• If priority is accuracy, then tabulation is better |
| **Line Diagram** | • Time Series is generally in x axis<br>• For wide fluctuation – log chart or ratio chart is used<br>• Two or more series of same unit – Multiple Line Chart<br>• Two or more series of different unit – Multiple Axis Chart |
| **Bar Diagram** | • Bar means rectangle of same width and of varying length drawn horizontally or vertically<br>• For comparable series – multiple or grouped bar diagrams can be used<br>• For data divided into multiple components – subdivided or component bar diagrams<br>• For relative comparison to whole, percentage bar diagrams or divided bar diagrams<br>• Vertical Bar Diagram: Useful for Data varying over Time and Quantitative Data<br>• Horizontal Bar Diagram: Useful for Data varying over Space and Qualitative Data |
| **Pie Chart** | • Used for circular presentation of relative data (% of whole)<br>• Summation of values of all components/segments are equated to 360 Degree (total angle of circle)<br>• **Segment angle =**<br><br>$$\frac{\text{(segment value x 360°)}}{\text{(total value)}}$$ |

**Statistical Description of Data – Frequency Distribution**

| | |
|---|---|
| **Frequency and Distribution** | • Frequency means number of times a particular observation is repeated.<br>• Frequency Distribution is table which contains observation or class intervals in one column and corresponding frequency in the other.<br>• Definition: A frequency distribution may be defined as a<br>  – tabular representation of statistical data, usually in an ascending order,<br>  – relating to a measurable characteristic<br>  – according to individual value or a group of values of the characteristic under study. |

| | | |
|---|---|---|
| **Types of Frequency Distribution** | **Ungrouped/ Simple Frequency Distribution** | • When there are limited number of distinct observations, frequency can be assigned to each one of them.<br>• This distribution is simple |
| | **Grouped Frequency Distribution** | • When there are large no. of observations, grouping is done among them (generally in ascending order).<br>• Each group is called as class interval and frequency is assigned to group and not individual values,<br>• this is called Grouped Frequency Distribution |

**Class Limit**

• For a class interval CL is the minimum and maximum value the class interval may contain
• Minimum Value – Lower Class Limit
• Maximum Value – Upper Class Limit

| Class Interval | Frequency | LCL | UCL |
|---|---|---|---|
| 10-19 | 10 | 10 | 19 |
| 20-29 | 5 | 20 | 29 |
| 30-39 | 8 | 30 | 39 |

**Classification of Grouped of Frequency Distribution**

Mutually Exclusive / Overlapping Classification

_20 exclude_

| Class | LCL | UCL |
|---|---|---|
| 10-20 | 10 | 20 |
| 20-30 | 20 | 30 |
| 30-40 | 30 | 40 |

• Here UCL an interval and LCL of next interval are same
• This is usually applicable for continuous variable.
• An observation which is equivalent to common class limit is excluded from the class interval where it is UCL and taken in the class where it is LCL.

Mutually Inclusive / Non-Overlapping Classification

| Class | LCL | UCL |
|---|---|---|
| 10-19 | 10 | 20 19 |
| 20-19 | 20 | 30 29 |
| 30-39 | 30 | 40 39 |

• There is no common class limits between two intervals.
• This is usually applicable to discrete variable.
• All observation including UCL and LCL will be taken in the same class interval as there is no confusion.

**Class Boundary**

In case of Exclusive / Overlapping Classification

Class Boundary = Class Limit

| Class | LCL | UCL | LCB | UCB |
|---|---|---|---|---|
| 10-20 | 10 | 20 | 10 | 20 |
| 20-30 | 20 | 30 | 20 | 30 |
| 30-40 | 30 | 40 | 30 | 40 |

In case of Inclusive / Overlapping Classification

Lower Class Boundary
LCB = LCL – 0.5
UCB = UCL + 0.5

| Class | LCL | UCL | LCB | UCB |
|---|---|---|---|---|
| 10-19 | 10 | 19 | 9.5 | 19.5 |
| 20-29 | 20 | 29 | 19.5 | 29.5 |
| 30-39 | 30 | 39 | 29.5 | 39.5 |

| Mid-Point / Class Mark / Mid Value of Class Interval | $\dfrac{LCL+UCL}{2}$ | | | $\dfrac{LCB+UCB}{2}$ | | |
|---|---|---|---|---|---|---|
| | • Useful in calculation of AM, GM, HM, SD in case of grouped frequency distribution | | | | | |

| Class Length/ Width or Size | UCB – LCB only |
|---|---|

| Cumulative Frequency | • Less than type: It shows no. of observations less than UCB<br>• More than type: It shows no. of observations more than UCB |
|---|---|

| Class Interval | Freq. | UCB | Less than type CF | More than type CF | Total of both CF |
|---|---|---|---|---|---|
| 44-48 | 3 | 48.5 | 3 | 33 | 36 |
| 49-53 | 4 | 53.5 | 7 | 29 | 36 |
| 54-58 | 5 | 58.5 | 12 | 24 | 36 |
| 59-63 | 7 | 63.5 | 19 | 17 | 36 |
| 64-68 | 9 | 68.5 | 28 | 8 | 36 |
| 69-73 | 8 | 73.5 | 36 | 0 | 36 |
| Total | 36 | | | | |

| Frequency Density | $\dfrac{\text{Class Frequency}}{\text{Class Length of class}}$ |
|---|---|
| Relative Frequency | $\dfrac{\text{Class frequency}}{\text{Total Frequency}}$<br>Its can have values between 0 and 1 |
| Percentage Frequency | $\dfrac{\text{Class frequency}}{\text{Total Frequency}} \times 100$ |
| Frequency Dist. Diagram – Histogram /area diag | • It is a convenient way to represent FD<br>• Comparison between frequency of two different classes possible<br>• It is useful to calculate mode also |
| Frequency Polygon | • Usually preferable for ungrouped frequency distribution<br>• Can be used for grouped also but only if class lengths are even |
| Ogives/ Cumulative Frequency | • This graph can be made by both type of Cumulative Frequency and called as Less than Ogive or More than Ogive<br>• It can be used for calculating quartiles, median |
| Frequency Curve | • It is a limiting form of Area Diagram (Histogram) or Frequency Polygon<br>• It is obtained by drawing smooth and free hand curve though the mid points<br>• Most used curve is Bell Shaped |

**Index Numbers**

| | |
|---|---|
| **Practical Examples of Index Numbers** | • Index numbers are convenient devices for **measuring relative changes (generally in %)** of differences from **time to time** or from **place to place** <br> • Series of numerical figures which show relative position <br> • Index Numbers show percentage changes rather than absolute amounts of change |
| **Data Selection** | • It **depends on the purpose** for which the index is used. <br> • Index numbers are often constructed from the **sample. Random sampling,** and if need be, a **stratified random sampling** can be used to ensure that sample is representative. <br> • Data should be **comparable** by ensuring consistency in selection method. |
| **Base Period** | • It is a **point of reference** in comparing various data. <br> • Standard point of comparison. <br> • The period should be **normal**. <br> • It should be **relatively recent** <br> • Choice of suitable base period is a temporary solution |
| **Use of Averages** | • The **geometric mean is better** in averaging relatives, <br> • But for most of the index's **arithmetic mean is used because of its simplicity** |
| **Price/ Quantity/ Value Relative** | For Individual Commodity, <br><br> **Current Period Price/ Quantity/ Value** <br> **Base Period Price/ Quantity/ Value** |
| **Link Relative** | $$\frac{P_1}{P_0}, \frac{P_2}{P_1}, \frac{P_3}{P_2}, ..., \frac{P_n}{P_{n-1}}$$ <br> Same can be created for quantities also |
| **Chain relatives** | When the above relatives are in respect to a fixed base period these are also called the chain relatives <br> $$\frac{P_1}{P_0}, \frac{P_2}{P_0}, \frac{P_3}{P_0}, ..., \frac{P_n}{P_0}$$ |
| **Formula for Chain Index (when direct data is not available)** | **Link relative of current year × Chain Index of previous year** <br> **100** <br><br> The chain index is an unnecessary complication unless of course where data for the whole period are not available or where commodity basket or the weights have to be changed. |
| **Limitations of Index Numbers** | • Chances of errors due to Sampling <br> • It gives broad trend not real picture <br> • Due to many methods, at times it creates confusion |
| **Usefulness of Index Numbers** | • Index numbers are very useful in deflating (eg. Nominal wages into real) <br> • Framing suitable policies in economics and business <br> • They reveal trends and tendencies in making important conclusions |

| | |
|---|---|
| | • They are used in time series analysis to study long-term trend, seasonal variations and cyclical developments |
| **Formula for Deflated Value** | $$\text{Deflated Value} = \frac{\text{Current Value}}{\text{Price Index of the current year}}$$ |
| **Shifted Price Index** | $$\frac{\text{Original Price Index}}{\text{Price Index of the year on which it has to be shifted}} \times 100$$ |
| **Unit Test** | • This test requires that the formula should be independent of the unit in which or for which prices and quantities are quoted.<br>• Except for the simple (unweighted) aggregative index all other formulae satisfy this test. |
| **Time Reversal Test** | • It is a test to determine whether a given method will work both ways in time, forward and backward.<br>• $P_{01} \times P_{10} = 1$<br>• Laspeyres' method and Paasche's method do not satisfy this test, but Fisher's Ideal Formula does. |
| **Factor Reversal Test** | • This holds when the product of price index and the quantity index should be equal to the corresponding value index.<br>• Symbolically<br>$$P_{01} \times Q_{01} = V_{01}$$<br>• Fisher's Index Number is ideal as it satisfies Unit, Time Reversal and Factor Reversal Test |
| **Circular Test** | • This property therefore enables us to adjust the index values from period to period without referring each time to the original base.<br>• It is an extension of time reversal test<br>• The test of this **shiftability of base** is called the circular test.<br>• This test is not met by Laspeyres, or Paasche's or the Fisher's ideal index.<br>• The weighted GM of relative, **simple geometric mean of price relatives** and the **weighted aggregative with fixed weights meet this test.**<br>(These methods are not in syllabus) |
| **Cost of Living Index (also called General Index)** | • CLI is defined as the **weighted AM of index numbers** of few groups of basic necessities.<br>• AM of group indices gives the General Index<br>• Generally, for calculating CLI; food, clothing, house rent, fuel & lightning and miscellaneous groups are taken into consideration.<br>• Examples of CLI: WPI, CPI, etc. |
| **Symbol** | • $P_{01}$ is the index for time 1 on 0<br>• $P_{10}$ is the index for time 0 on 1 |

$P_{01}$

## Measures of Central Tendency

### Arithmetic Mean

| Property 1 | If all the observations are constant, AM is also constant |
|---|---|
| Property 2 | the algebraic sum of deviations of a set of observations from their AM is zero |
| Property 3 | AM is affected both due to change of origin and scale<br>If $y = a + bx$ then $\bar{y} = a + b\bar{x}$ |
| Property 4 | Combined AM<br><br>$$\bar{x}_c = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$ |
| General Review | • AM is best measure of central tendency<br>• AM is based on all observations<br>• AM is affected by sampling fluctuations<br>• AM is amenable to mathematical property<br>• AM cannot be used in case of open end classification |

### Median

| Property 1 | For a set of observations, the sum of absolute deviations is minimum, when the deviations are taken from the median.<br><br>$$\sum |x_i - Me| \quad \text{modulus}$$ |
|---|---|
| Property 2 | Median is also affected by both change of origin and scale. |
| General Review | • Median is also called as positional average<br>• Median is not based on all observations<br>• Median is not affected by sampling fluctuations<br>• Median is best measure of central tendency in case of open end classification |

### Partition Values

| Meaning | • These may be defined as **values dividing** a given **set of observations** into number of **equal parts**<br>• When we want to divide the given set of observations into two equal parts, we consider median, similarly there are quartiles, deciles, percentiles |
|---|---|

| Name of PV | No. of equal parts | No. of PVs | Symbol |
|---|---|---|---|
| Median | 2 | 1 | **Me** |
| Quartile | 4 | 3 | $Q_1, Q_2, Q_3$ |

| | Decile | 10 | 9 | $D_1, D_2, ..., D_9$ |
| --- | --- | --- | --- | --- |
| | Percentile | 100 | 99 | $P_1, P_2, ..., P_{99}$ |

## Mode – Concept/ Formula

| Meaning | Mode is the **value** that **occurs the maximum** number of times |
| --- | --- |
| **Special Thing about Mode** | • If two or more observations are having maximum frequency then there are **multiple modes** [multimodal distribution] <br> • If there are **exactly two** modes then distribution is called as **Bimodal** Distribution <br> • If all observations are having same frequency then distribution has **no mode** <br> • We can say that Mode is **not rigidly defined** |
| **Property 1** | If all the observations are constant, mode is also constant |
| **Property 2** | Mode is also affected both due to change of origin and scale |
| **General Review** | • Mode is not based on all observations <br> • Mode is not rigidly defined <br> • Mode is not amenable to Mathematical Property |

## Relationship between Mean, Median and Mode

| **In case of Symmetric Distribution** | **Mean = Median = Mode** |
| --- | --- |
| **In case of Moderately Skewed Distribution (Empirical relationship)** | **Mean – Mode = 3 (Mean – Median)** |

## Geometric Mean

| **Definition** | For a given set of $n$ **positive observations**, the geometric mean is defined as the $n^{th}$ root of the product of the observations |
| --- | --- |
| **Property 1** | Logarithm of G for a set of observations is the AM of the logarithm of the observations <br><br> $$\log G = \frac{1}{n} \sum \log x$$ |
| **Property 2** | If all the observations are constant, GM is also constant |
| **Property 3** | $\text{GM of } z = \text{GM of } x \times \text{GM of } y$ |
| **Property 4** | $\text{GM of } z = \dfrac{\text{GM of } x}{\text{GM of } y}$ |

## Harmonic Mean

| | |
|---|---|
| **Definition** | For a given set of **non-zero** observations, harmonic mean is defined as the **reciprocal of the AM of the reciprocals of the observation** |
| **Property 1** | If all observations are constant HM is also constant |

## Use of GM and HM

| | |
|---|---|
| Both | Both are used for calculating average rates |
| GM | Appropriate for rates having percentages |
| HM | Appropriate for rates other than percentages |

## Measures of Dispersion

| | | |
|---|---|---|
| **Meaning of Measure of Dispersion** | • Dispersion for a given set of observations may be defined as<br>• the **amount of deviation** of the observations,<br>• usually, from an **appropriate** measure of **central tendency** | |
| **Types of Measure of Dispersion** | **Absolute Measures of Dispersion** | • These are with units<br>• These are not useful for comparison of two variables with different units.<br>• Example: Range, Mean Deviation, Standard Deviation, Quartile Deviation |
| | **Relative Measures of Dispersion** | • These are unit free measures<br>• These are useful for comparison of two variables with different units.<br>• Example: Coefficient of Range, Coefficient of Mean Deviation, Coefficient of variation, Coefficient of Quartile Deviation |

## Range

| | |
|---|---|
| **Property 1** | • **Not affected** by change of **origin**<br>• Affected by **change of scale (only value)**<br>• **No impact of sign** of change of scale<br>• Note: **Measure of Dispersion can never be negative** |
| **General Review** | • Not Based on All Observations<br>• Easy to Compute |

## Mean Deviation

| | |
|---|---|
| **Meaning** | • Mean deviation is defined as the<br>• **arithmetic mean** of the<br>• **absolute deviations** of the observations<br>• from an **appropriate measure** of central tendency |
| **Property 1** | Mean Deviation takes its **minimum value** when deviations are taken from **Median** |

| Property 2 | Change of Origin – **No Affect**, Change of Scale – **Affect of value not sign** |
|---|---|
| **General Review** | <ul><li>Based on **all observations**</li><li>Improvement over Range</li><li>**Difficult to compute**</li><li>**Not amenable to Mathematical Property** because of usage of **Modulus**</li></ul> |

## Standard Deviation

| **Meaning** | <ul><li>Improvement over Mean Deviation</li><li>It is defined as the **root mean square deviation** when the deviations are <u>taken from the AM</u> of the observations</li></ul> |
|---|---|
| **Coefficient of Variation** | $\dfrac{SD_x}{\overline{x}} \times 100$ |
| **SD for any two numbers** | $SD = \dfrac{\lvert a-b \rvert}{2}$ |
| **SD for first n natural numbers** | $s = \sqrt{\dfrac{n^2-1}{12}}$ |
| **Property 1** | If all the observations are constant, SD is **ZERO** |
| **Property 2** | No effect of change of origin but affected by change of scale in the magnitude (ignore sign) |
| **Property 3** | $SD_c = \sqrt{\dfrac{n_1 s_1^2 + n_2 s_2^2 + n_1 d_1^2 + n_2 d_2^2}{n_1 + n_2}}$ <br> $d_1 = \overline{x}_c - \overline{x}_1$ <br> $d_2 = \overline{x}_c - \overline{x}_2$ |

## Quartile Deviation

| **Meaning** | It is semi-inter quartile range |
|---|---|
| **General Review** | <ul><li>It is the **best measure** of dispersion for **open-end** classification</li><li>It is also **less affected** due to sampling fluctuations</li><li>Like other measures of Dispersion, **QD** is also not affected by change of origin but affected by scale ignoring sign</li></ul> |

## Correlation and Regression

### Bivariate Data

| | |
|---|---|
| **Definition** | • When data are collected on two variables **simultaneously**, they are known as **bivariate data** <br> • and the corresponding frequency distribution, derived from it, is known as **Bivariate Frequency Distribution** |
| **Marginal Distribution** | • It is the frequency distribution of **one variable** (x or y) across the other variable's **full range of values** <br> • **Number of Marginal Distribution = 2** |
| **Conditional Distribution** | • It is the frequency distribution of **one variable** (x or y) across a **particular sub-population** of the other variable. <br> • **No. of Conditional Distributions = m + n** <br> *m = no. of class interval of x* <br> *n = no. of class interval of y* |

### Scatter Diagram

| | |
|---|---|
| **Concept Points** | • It helps us to find **Nature** and **Relative Strength** of Correlation <br> • It is useful for **Non-Linear** Correlation also <br> • It **cannot** be used to determine **value** <br> • Diagrams are **time taking** |

### Karl Pearson's Correlation Coefficient

| | | |
|---|---|---|
| **How to Calculate** | Correlation Coefficient is the ratio of covariance with product of standard deviations | |
| **Property 1** | The Coefficient of Correlation is a **unit-free measure** | |
| **Property 2** | Value lies from **-1 to +1** | |
| **Property 3** | **Change of Origin** | No impact |
| | **Change of Scale** | No impact of value, but if change of scale of both variables are of **different sign** then **sign of r** will also change |

| | | |
|---|---|---|
| **Interpretation of Value of r** | **Value of r** | **Interpretation** |
| | -1 | Perfect Negative |
| | Between -1 and 0 | Negative |
| | Closer to -1 | Strong Negative |
| | Far from -1 | Weak Negative |
| | 0 | No Correlation |
| | Between 0 and 1 | Positive |
| | Far from +1 | Weak Positive |
| | Near to +1 | Strong Positive |
| | +1 | Perfect Positive |

## Spearman's Rank Correlation Coefficient

| Usage | • find the level of **agreement (or disagreement)** between two judges so far as assessing a **qualitative characteristic (attribute)** is concerned<br>• Use in case of ranks |
|---|---|
| **Ranking in case of Tie** | In case of tie, simple average of ranking should be assigned to tied values |

## Coefficient of Concurrent Deviations

| Usage | A very **quick, simple** and **casual** method of finding correlation when we are not serious about the magnitude of the two variables |
|---|---|

## Regression Basics

| Meaning | Estimation of one variable for a **given value** of another variable on the basis of an **average mathematical relationship** between the two variables | |
|---|---|---|
| **Requirements** | • Estimation of Y when X is given<br>• Estimation of X when Y is given | |
| **General Points** | **Perfect Correlation** | • When linear relationship exists between two variables, correlation is perfect.<br>• Perfect Correlation is represented by a linear equation and this equation can be used for regression purpose directly.<br>• Same equation can be used in both ways |
| | **Imperfect Correlation** | • In case of imperfect correlation there is no definite line and equation<br>• We will use method of least square to estimate both regression lines |
| **Formula of Regression Equations/ Lines** | Estimation of Y when X is given | • Use Regression line of **Y on X**<br>• Equation Format:<br>$$Y - \overline{Y} = b_{yx}(X - \overline{X})$$<br>$b_{yx}$ is regression coefficient of Y on X |
| | Estimation of X when Y is given | • Use Regression line of **X on Y**<br>• Equation Format:<br>$$X - \overline{X} = b_{xy}(Y - \overline{Y})$$<br>$b_{xy}$ is regression coefficient of X on Y |
| **Property 1** | Change of Origin and Scale<br>• Origin: No Impact<br>• Scale: If original pair is $x, y$ and modified pair is $u, v$<br>$$b_{vu} = b_{yx} \times \frac{\text{change of scale of } y}{\text{change of scale of } x}$$ | |

| | |
|---|---|
| | $$b_{uv} = b_{xy} \times \dfrac{\text{change of scale of x}}{\text{change of scale of y}}$$ |
| **Property 2** | Two regression lines (if not identical) will intersect at the point [means] $(\bar{x}, \bar{y})$ |
| **Property 3** | Relation between Correlation and Regression Coefficients $$r_{xy} = \pm\sqrt{b_{xy} \times b_{yx}}$$ $r_{xy}, b_{xy}, b_{yx}$ will always have same sign |

**Probable Error**

| | |
|---|---|
| **Use** | • Correlation is calculated using sample, value for sample may differ from population, this difference is probable error<br>• If there is significant probable error, there is no evidence of real correlation |
| **Limits of Sample Correlation Coefficient** | $r \pm PE$ |

| **How to check evidence of Correlation using PE** | Case | Conclusion |
|---|---|---|
| | If r is less than PE | There is no evidence of correlation |
| | If r is greater than six times of PE | The presence of correlation is certain |
| | Since r lies from -1 to +1 | PE can never be negative |

**Coefficient of Determination and Non-Determination**

| | |
|---|---|
| **Coefficient of Determination**<br>Accounted Variance/ Explained Variance | $r^2$ |
| **Coefficient of Non-Determination**<br>Unaccounted Variance/ Unexplained Variance | $1 - r^2$ |